



# **Deep Learning-Based Environmental Forecasting: Enhancing Accuracy with Bidirectional LSTMs**

---

**By Anonymous 3.0 - Data\_Crunch\_071  
General Sir John Kotelawala Defence University**

---

# 1. Problem Understanding & Dataset Analysis

## 1.1. Forecasting Objective

The objective of this study is to develop an accurate environmental forecasting model that can predict multiple meteorological variables, including,

- Average Temperature
- Solar Radiation
- Rainfall Amount
- Wind Speed
- Wind Direction

By leveraging historical weather patterns, this model aims to generate reliable future forecasts, which can be applied to domains such as agriculture, energy management, disaster preparedness, and water resource planning.

## 1.2. Key Findings from Data Analysis

### Understanding the Dataset

The dataset consists of daily time series records spanning multiple years. Each record contains various meteorological measurements, including temperature, radiation, rainfall, wind speed, and wind direction.

### Challenges Identified in the Data

During initial data analysis, several issues were detected,

1. Date Formatting Issues
  - The dataset contains date inconsistencies, especially in leap years (e.g., February 29).
  - Some entries had incorrect timestamps, requiring reformatting.
2. Outliers & Invalid Measurements
  - Temperature values were recorded in both Kelvin and Celsius, leading to inconsistencies.
3. Seasonal Patterns
  - Strong seasonality was observed in environmental variables, indicating the need for seasonal features to improve forecasting accuracy.

## Preprocessing Steps & Justification

To ensure high-quality input data, the following data preprocessing techniques were applied,

1. Date Handling
  - The year, month, and day columns were combined into a proper datetime format while addressing leap year issues.
2. Missing Value Handling
  - A forward-fill and backward-fill approach was applied to preserve time series continuity.
3. Temperature Standardization
  - Any temperature values above 100°C were converted from Kelvin to Celsius.
4. Physical Constraints Enforcement
  - Negative rainfall and evapotranspiration values were set to zero.
  - Negative wind speeds were corrected using absolute values.
  - Wind direction values were restricted to a valid range of 0-360°.
5. Outlier Detection & Treatment
  - The Interquartile Range (IQR) method was used with a 3× multiplier to cap extreme values instead of removing them, preserving useful data.

## 2. Feature Engineering & Data Preparation

### 2.1. Feature Creation Techniques

To improve predictive performance, new features were created based on domain knowledge and data analysis,

1. Temporal Features
  - Basic Features: Extracted month, day, and day of the year.
  - Cyclical Encoding: Converted seasonal features using sine and cosine transformations to capture periodicity.
2. Lag Features
  - Created lagged variables (1, 3, 7, 14 days) to capture short- and medium-term dependencies.
3. Rolling Statistics
  - Moving averages and rolling standard deviations over different window sizes (3, 7, and 14 days) to capture trends and variability.
4. Interaction Terms
  - Temperature-Rainfall Interaction: Models the effect of temperature on precipitation.
  - Wind-Rainfall Interaction: Helps understand how wind affects rainfall distribution.

### 2.2. Feature Selection & Data Transformations

1. Feature Selection: Selected the top 15 most correlated features for each target variable.
2. Normalization: RobustScaler for input features to handle outliers effectively  
MinMaxScaler for target variables to constrain outputs to a reasonable range (0-1).
3. Sequence creation: 90-day sequences were created using a sliding window approach This sequence length was selected to capture both short-term patterns and seasonal effects.
4. Train-validation split: 80% training, 20% validation chronological split preserving time series integrity This approach respects the temporal nature of the data without data leakage.

## 3. Model Selection & Justification

### 3.1. Evaluated Models

Several models were tested for forecasting:

1. **ARIMA:** Good for univariate time series but struggles with multiple variables.
2. **XGBoost:** Works well with structured data but does not handle temporal dependencies effectively.
3. **Prophet:** Designed for seasonal data but lacks deep learning capabilities.
4. **LSTM (Long Short-Term Memory):** Chosen due to its ability to learn sequential patterns.

### 3.2. Justification for Bidirectional LSTM

- Captures past and future dependencies better than unidirectional models.
- Models complex environmental relationships between multiple meteorological variables.
- Uses hierarchical feature extraction through stacked layers.

### Model Architecture

- Two Bidirectional LSTM layers for deep sequence learning.
- Unidirectional LSTM layer for additional feature extraction.
- Dropout Regularization (0.2-0.3) to prevent overfitting.
- Dense output layer for final predictions.

### 3.3. Hyperparameter Optimization

- **Sequence Length:** 90 days (optimized for seasonal trends).
- **Layer Configuration:** 128 → 96 → 64 → 48 neurons.
- **Batch Size:** 32 (optimized for stable training).
- **Learning Rate:** 0.001, adjusted dynamically using ReduceLROnPlateau.
- **Early Stopping:** Stops training if validation loss does not improve for 10 epochs.

### 3.4. Validation Approach

- Chronological Train-Test Split (80/20%) to maintain time integrity.
- Walk-forward validation for realistic forecasting performance evaluation.

## 4. Performance Evaluation & Error Analysis

### 4.1. Evaluation Metrics

1. SMAPE (Symmetric Mean Absolute Percentage Error)
  - Measures relative accuracy, making it useful for comparing different environmental variables.
  - Less sensitive to small errors but may underestimate minor fluctuations.
2. RMSE (Root Mean Squared Error)
  - Highlights large prediction errors, making it useful for detecting extreme weather deviations.
  - Sensitive to outliers, which can inflate the error if extreme events occur.
3. MAE (Mean Absolute Error)
  - Provides an intuitive measure of the average error magnitude.
  - Does not emphasize large errors as much as RMSE, making it useful for general forecasting accuracy.

### 4.2. Residual Analysis & Model Limitations

1. Systematic Prediction Errors for Extreme Weather Events
  - The model struggled with sudden spikes in rainfall, wind speed, and temperature.
  - Extreme weather events were often underestimated or missed due to their rarity in the dataset.
2. Error Accumulation in Long-Term Forecasts
  - Small prediction errors compounded over time, reducing accuracy for forecasts beyond 30 days.
  - The model had difficulty capturing seasonal changes beyond its training window.
3. High Computational Cost of Bidirectional LSTM Networks
  - Training was time-intensive due to the sequential processing nature of LSTMs.
  - The bidirectional architecture doubled the computational load, making it expensive for real-time applications.

## 5. Interpretability & Business Insights

### 5.1. Applications of the Forecasting Model

1. Agriculture
  - Accurate weather forecasts help farmers plan irrigation, fertilization, and harvesting.
  - Predicting rainfall and temperature trends enables better pest and disease management.
  - Helps optimize crop selection based on expected seasonal conditions.
2. Energy Sector
  - Solar power generation depends on accurate sunlight predictions.
  - Wind speed forecasts improve wind energy management and prevent turbine damage from sudden gusts.
  - Helps energy companies balance supply and demand by forecasting weather-driven power consumption changes.
3. Disaster Preparedness
  - Timely weather predictions help authorities issue early warnings for storms, floods, and heatwaves.
  - Reduces property damage and loss of life by aiding in evacuations and emergency response planning.
  - Helps insurance companies assess weather-related risks for businesses and homeowners.
4. Water Resource Management
  - Predicts drought conditions, helping water agencies allocate resources efficiently.
  - Supports reservoir management by forecasting precipitation and runoff patterns.
  - Helps in flood control planning by predicting extreme rainfall events.

### 5.2. Recommendations for Improvement

#### Model Enhancements

1. Combine LSTM with Ensemble Learning
  - Using LSTM with traditional models (like XGBoost or Random Forests) can improve overall accuracy.
  - This approach helps capture both long-term trends (LSTM) and short-term variations (tree-based models).



## 2. Introduce Attention Mechanisms

- Attention mechanisms prioritize important time steps, improving long-term forecasting accuracy.
- This is useful when the model needs to focus on key weather patterns from past data.

## Feature Engineering Improvements

### 1. Incorporate External Climate Indices (ENSO, NAO, etc.)

- Climate indices like El Niño-Southern Oscillation (ENSO) and North Atlantic Oscillation (NAO) influence weather patterns.
- Including these indices can improve the model's ability to predict seasonal variations.

### 2. Utilize Spatial Meteorological Data

- Weather conditions are influenced by geographic and topographic factors.
- Incorporating satellite data, elevation maps, and regional weather station data can improve model performance.

## 6. Innovation & Technical Depth

### 6.1. Key Innovations

1. Use of Bidirectional LSTMs for Enhanced Pattern Recognition
  - Traditional LSTM (Long Short-Term Memory) networks process data only in one direction (past to future).
  - Bidirectional LSTMs process data in both forward and backward directions, improving the model's ability to capture complex weather patterns.
  - This is especially useful for detecting seasonal trends, sudden weather shifts, and long-term dependencies.
2. Feature Engineering Including Rolling Statistics and Interaction Terms
  - Rolling statistics (e.g., moving averages, standard deviation) help smooth out short-term fluctuations while retaining key trends.
  - Interaction terms capture relationships between different variables (e.g., combining temperature and humidity to improve rainfall predictions).
  - These engineered features enhance model accuracy by providing additional meaningful insights.
3. Hyperparameter Optimization Using Dynamic Learning Rate Adjustments
  - Instead of using a fixed learning rate, the model dynamically adjusts it based on performance.
  - This helps the model learn faster initially, then fine-tune itself more precisely as training progresses.
  - Prevents overfitting and ensures better convergence, leading to higher forecasting accuracy.
4. Walk-Forward Validation for Time-Series Robustness
  - Unlike traditional validation methods, walk-forward validation ensures the model is tested on never-before-seen future data.
  - This method better simulates real-world forecasting conditions, making the model more reliable.
  - It continuously updates the training dataset as new data becomes available, improving adaptability over time.

## 7. Conclusion

### 7.1. Summary of Findings

1. Strong Forecasting Accuracy with Bidirectional LSTM
  - The Bidirectional LSTM model successfully captured complex weather patterns by analyzing past and future dependencies.
  - It outperformed traditional forecasting models in detecting seasonal trends and fluctuations in environmental variables.
2. Feature Engineering Significantly Improved Performance
  - The inclusion of rolling statistics, interaction terms, and external climate indices led to better predictions.
  - These engineered features allowed the model to detect subtle relationships between variables like temperature, humidity, and rainfall.
3. High-Quality Data Preprocessing Pipeline
  - Proper data cleaning, normalization, and outlier handling ensured that the model was trained on reliable and accurate data.
  - The walk-forward validation approach enhanced model robustness by testing it on unseen future data.

### 7.2. Challenges & Future Improvements

1. Handling Extreme Weather Variability
  - The model struggled to predict sudden extreme events like storms and heatwaves due to their rarity in historical data.
  - Future improvements could involve data augmentation techniques or external weather sources to improve extreme event detection.
2. Reducing Computational Costs for Real-Time Deployment
  - Training and deploying Bidirectional LSTMs require high computational power, making real-time forecasting challenging.
  - Optimizations such as model compression, GPU acceleration, or switching to lightweight architectures (e.g., GRUs) could enhance efficiency.
3. Exploring Hybrid Models for Improved Forecasting
  - Combining deep learning (LSTMs) with statistical models (ARIMA, XGBoost, etc.) may enhance long-term forecasting accuracy.
  - Hybrid approaches could help balance short-term precision with long-term trend prediction.

